

Cuestionario de teoría 1

Antonio Álvarez Caballero
analca3@correo.ugr.es

3 de abril de 2016

1. Cuestiones

Cuestión 1. Identificar, para cada una de las siguiente tareas, qué tipo de aprendizaje automático es el adecuado (supervisado, no supervisado, por refuerzo) y los datos de aprendizaje que deberíamos usar. Si una tarea se ajusta a más de un tipo, explicar cómo y describir los datos para cada tipo.

- a) Categorizar un grupo de animales vertebrados en pájaros, mamíferos, reptiles, aves y anfibios.
- b) Clasificación automática de cartas por distrito postal.
- c) Decidir si un determinado índice del mercado de valores subirá o bajará dentro de un periodo de tiempo fijado.

Solución. Los tipos de aprendizaje son:

- a) Aprendizaje supervisado. Lo importante aquí es el conjunto de datos, ya que cada uno de estas clases tiene características bastante bien definidas (pelo, plumas, número de patas, número de alas, piel seca o húmeda, escamas, órganos respiratorios, ...). Por esta misma razón estaría bien intentar realizar una aproximación por diseño, aunque habría que realizar primero un análisis de las variables y ver existe alguna pequeña superposición entre clases.
- b) Aprendizaje no supervisado. Aquí no hace falta que se sepa realmente la etiqueta de los datos, simplemente con que el sistema sea capaz de agrupar las cartas que pertenezcan al mismo distrito postal (que tengan el mismo número manuscrito).
- c) Por refuerzo. No siempre podemos basarnos en datos que ya tenemos para decidir esto (los datos de hace 4 años no tienen por qué ser relevantes ahora), por lo que una aproximación por refuerzo irá aprendiendo sobre la marcha con la actividad actual del mercado de valores.

Cuestión 2. ¿Cuáles de los siguientes problemas son más adecuados para una aproximación por aprendizaje y cuáles más adecuados para una aproximación por diseño? Justificar la decisión.

- a) Determinar el ciclo óptimo para las luces de los semáforos en un cruce con mucho tráfico.
- b) Determinar los ingresos medios de una persona a partir de sus datos de nivel de educación, edad, experiencia y estatus social.
- c) Determinar si se debe aplicar una campaña de vacunación contra una enfermedad.

Solución. Las aproximaciones son:

- a) Aproximación por diseño. Este es un problema que puede llegar a ser crítico si se produce algún fallo, luego una aproximación por diseño (exhaustivo) dará mejor resultado en una primera instancia que una aproximación por aprendizaje.
- b) Aproximación por aprendizaje. Los ingresos medios de una persona no están directamente relacionados con los datos proporcionados, luego deberíamos intentar aprender de esos datos y realizar algún modelo que nos permita *predecir* este dato de forma aceptable.
- c) Aproximación por aprendizaje. Se podría realizar aprendizaje sobre algún conjunto de datos que incluyera (por ejemplo) un factor de contagio de la enfermedad, probabilidad de secuelas críticas, etc... para, dada una enfermedad, poder decidir si es necesaria o no dicha campaña.

Cuestión 3. Construir un problema de *aprendizaje desde datos* para un problema de selección de fruta en una explotación agraria (ver transparencias de clase). Identificar y describir cada uno de sus elementos formales. Justificar las decisiones.

Solución. Determinemos los aspectos clave:

- \mathcal{X} : Un espacio n-dimensional con las características de los mangos. Podríamos tomar un espacio 3-dimensional donde las características son (color, tamaño, textura).
- \mathcal{Y} : Podemos usar clasificación binaria: 1 si el mango es bueno y -1 si no lo es.
- \mathcal{A} : Se propone aprendizaje supervisado, ya que ...
- Conjunto de entrenamiento: Un análisis previo de la cosecha midiendo cada una de las características propuestas y determinando si el mango es bueno o no.
- \mathcal{H} : En principio podríamos intentar buscar en las funciones lineales, ya que son las más sencillas. Si no consiguiéramos un buen clasificador, ya iríamos subiendo la complejidad de estas funciones en busca de un clasificador aceptable.

Cuestión 4. Suponga un modelo PLA y un dato $x(t)$ mal clasificado respecto de dicho modelo. Probar que la regla de adaptación de pesos del PLA es un movimiento en la dirección correcta para clasificar bien $x(t)$.

Solución. Definimos primero el término $y(t)w^T(t)x(t)$ como, informalmente, el resultado de clasificar $x(t)$ con $w(t)$. Claramente es > 0 si está bien clasificado y < 0 si no lo está, ya que si está bien clasificado tenemos que $y(t) = \text{sign}(w^T(t)x(t))$, por lo cual su producto es positivo, y en caso contrario $y(t) \neq \text{sign}(w^T(t)x(t))$, lo cual su producto es negativo.

Usando esto, veamos que $y(t)w^T(t+1)x(t) > y(t)w^T(t)x(t)$. Es claro ya que si expandimos el primer miembro con la regla de actualización del PLA vemos que:

$$y(t)w^T(t+1)x(t) = y(t)(w(t) + y(t)x(t))^T x(t) = y(t)w^T(t)x(t) + \underbrace{\|x(t)\|^2}_{>0}$$

Por lo cual la desigualdad de arriba queda clara. Con esto deducimos que la clasificación de $x(t+1)$ siempre será *mejor* que la de $x(t)$, ya que el valor de su clasificación es mayor, lo cual hace que pueda superar a 0, implicando que el dato estaría bien clasificado. Podría no superar a 0, pero en tal caso estaría más cerca de estar bien clasificado.

Cuestión 5. Considere el enunciado del ejercicio 2 de la sección FACTIBILIDAD DEL APRENDIZAJE de la relación de apoyo.

- a) Si $p = 0,9$, ¿cuál es la probabilidad de que S produzca una hipótesis mejor que C ?

- b) ¿Existe un valor de p para el cual es más probable que C produzca una hipótesis mejor que S ?

Solución. Las soluciones son:

- a) Queremos determinar cuál es la probabilidad de que el error fuera de la muestra sea menor para S que para C :

$$\begin{aligned} P[E_{out}(S(\mathcal{D})) < E_{out}(C(\mathcal{D}))] &= P[E_{out}(h_1) < E_{out}(h_2)] \\ &= P[P[f(x) \neq h_1] < P[f(x) \neq h_2]] \\ &= P[P[f(x) = -1] < P[f(x) = +1]] \\ &= P[1 - p < p] = P[0,1 < 0,9] = 1 \end{aligned}$$

- b) En este caso razonamos igual que en el anterior, pero dándole la vuelta a la desigualdad y sin determinar p :

$$\begin{aligned} P[E_{out}(S(\mathcal{D})) > E_{out}(C(\mathcal{D}))] &= P[E_{out}(h_1) > E_{out}(h_2)] \\ &= P[P[f(x) \neq h_1] > P[f(x) \neq h_2]] \\ &= P[P[f(x) = -1] > P[f(x) = +1]] \\ &= P[1 - p > p] = P[0,5 > p] = 1 \end{aligned}$$

Es claro que dicha probabilidad es 1 si y solamente si $p < 0,5$

Cuestión 6. La desigualdad de Hoeffding modificada nos da una forma de caracterizar el error de generalización con una cota probabilística

$$\mathbb{P}[|E_{out}(g) - E_{in}(g)| > \epsilon] \leq 2Me^{-2N\epsilon^2}$$

para cualquier $\epsilon > 0$. Si fijamos $\epsilon = 0,05$ y queremos que la cota probabilística $2Me^{-2N\epsilon^2}$ sea como máximo 0,03, ¿cuál será el valor más pequeño de N que verifique estas condiciones si $M = 1$? Repetir para $M = 10$ y para $M = 100$.

Solución. Sólo tenemos que despejar N de la ecuación para conseguir los valores deseados. Como queremos que la cota sea 0,03, igualamos y despejamos:

$$2Me^{-2N\epsilon^2} = 0,03$$

$$e^{-2N\epsilon^2} = \frac{0,015}{M}$$

$$-2N\epsilon^2 = \log\left(\frac{0,015}{M}\right)$$

$$N = \lceil -\log\left(\frac{0,015}{M}\right) \frac{1}{2\epsilon^2} \rceil$$

Ahora sólo tenemos que introducir $\epsilon = 0,05$ y los distintos valores de M .

$$M = 1 \rightarrow N = 840$$

$$M = 10 \rightarrow N = 1301$$

$$M = 100 \rightarrow N = 1761$$

$$M = 1000 \rightarrow N = 2222$$

Cuestión 7. Consideremos el modelo de aprendizaje M -intervalos donde $h : \mathbb{R} \rightarrow \{-1, +1\}$ y $h(x) = +1$ si el punto está dentro de cualquiera de m intervalos arbitrariamente elegidos y $h(x) = -1$ en otro caso. ¿Cuál es el más pequeño punto de ruptura para este conjunto de hipótesis?

Solución. El modelo de M -intervalos es una generalización del modelo del *intervalo*. Sabemos que para este modelo el punto de ruptura es $k = 3$. Veamos entonces si conseguimos alguna función lineal $aM + b$; $a, b \in \mathbb{N}$ que nos dé dicho valor de k para el caso $M = 1$. Claramente los posibles valores son $a = 2 \wedge b = 1$ o $a = 1 \wedge b = 2$. Probamos con el primero, dándonos la función $2M + 1$, viendo si el modelo es capaz de clasificar completamente $2M$ dicotomías y no puede con $2M + 1$.

Claramente el modelo puede separar correctamente un conjunto de $2M$ puntos. Supongamos que tenemos M 1's y M -1's. Recubriendo cada uno con un intervalo (sin alcanzar ningún -1) podemos clasificar cualquier dicotomía de este conjunto. Al haber separado un conjunto, podemos decir que NO es punto de ruptura.

Veamos para un conjunto con $2M + 1$ puntos: la condición que necesitamos es que NO pueda clasificar todas las dicotomías de ningún conjunto. Es claro que no se puede: tomamos una dicotomía con M -1's y $M + 1$ 1's. Si estos puntos están alternados $(1, -1, 1, \dots, 1, -1, 1)$, no podemos clasificarlos con M intervalos.

Cuestión 8. Suponga un conjunto de k^* puntos x_1, x_2, \dots, x_{k^*} sobre los cuales la clase H implementa $< 2^{k^*}$ dicotomías. ¿Cuáles de las siguientes afirmaciones son correctas?

- a) k^* es un punto de ruptura
- b) k^* no es un punto de ruptura
- c) todos los puntos de ruptura son estrictamente mayores que k^*
- d) todos los puntos de ruptura son menores o iguales a k^*
- e) no conocemos nada acerca del punto de ruptura

Solución. La respuesta correcta es la e). Para llegar a una conclusión tendría que ser *para todo conjunto de k^* puntos*.

Cuestión 9. Para todo conjunto de k^* puntos, H implementa $< 2^{k^*}$ dicotomías. ¿Cuáles de las siguientes afirmaciones son correctas?

- a) k^* es un punto de ruptura
- b) k^* no es un punto de ruptura
- c) todos los $k \geq k^*$ son puntos de ruptura
- d) todos los $k < k^*$ son puntos de ruptura
- e) no conocemos nada acerca del punto de ruptura

Solución. Las respuestas correctas son a) y c). Como son todos los conjuntos, 2^{k^*} acota las dicotomías que implementa H . Además, cualquier valor mayor a dicho k^* también lo acotará.

Cuestión 10. Si queremos mostrar que k^* es un punto de ruptura, ¿cuáles de las siguientes afirmaciones nos servirían para ello?:

- a) Mostrar que existe un conjunto de k^* puntos x_1, \dots, x_{k^*} que H puede separar (*shatter*).
- b) Mostrar que H puede separar cualquier conjunto de k^* puntos.

- c) Mostrar un conjunto de k^* puntos x_1, \dots, x_{k^*} que H no puede separar.
- d) Mostrar que H no puede separar ningún conjunto de k^* puntos.
- e) Mostrar que $m_H(k) = 2^{k^*}$

Solución. La solución es d). Si separara alguno, estaríamos en el caso de la cuestión 8. Si los separara todos, en el caso de la 9.

Cuestión 11. Para un conjunto H con $d_{VC} = 10$, ¿qué tamaño muestral se necesita (según la cota de generalización) para tener un 95 % de confianza de que el error de generalización sea como mucho 0,05?

Solución. Tal y como se explica en las transparencias, se resolverá la cota del tamaño de la muestra iterativamente sobre la fórmula apropiada.

```

1  #!/usr/bin/env python
2
3  import math
4
5  EPSILON = 0.05
6  DELTA = 0.05
7  dVC = 10
8
9  """
10 Aplicamos la formula explicada en teoria
11 """
12 def newN(n):
13     result = 8 / (EPSILON**2)
14     result *= math.log(4*((2*n)**dVC + 1)/DELTA)
15     return result
16
17 if __name__ == "__main__":
18
19     N = 10000
20     new_N = newN(N)
21
22     # Iteramos hasta que la diferencia entre iteraciones sea baja
23     while math.fabs(new_N - N) > 1e-6:
24         N = new_N
25         new_N = newN(N)
26
27     # Imprimimos el entero por encima del resultado
28     print(math.ceil(N))

```

El tamaño resultado de ejecutar este script es de 452957.

Cuestión 12. Consideremos un escenario de aprendizaje simple. Supongamos que la dimensión de entrada es uno. Supongamos que la variable de entrada x está uniformemente distribuida en el intervalo $[-1, 1]$ y el conjunto de datos consiste en 2 puntos $\{x_1, x_2\}$ y que la función objetivo es $f(x) = x^2$. Por tanto el conjunto de datos completo es $\mathcal{D} = \{(x_1, x_1^2), (x_2, x_2^2)\}$. El algoritmo de aprendizaje devuelve la línea que ajusta estos dos puntos como g (i.e. H consiste en funciones de la forma $h(x) = ax + b$).

- a) Dar una expresión analítica para la función promedio $\bar{g}(x)$.
- b) Calcular analíticamente los valores de E_{out} , *bias* y *var*.

Solución. Primero definamos explícitamente $g^{(\mathcal{D})}(x)$. Como es la recta que corta a los puntos de \mathcal{D} , se define en forma punto pendiente como $g^{(\mathcal{D})}(x) = \frac{x_2^2 - x_1^2}{x_2 - x_1}(x - x_1) + x_1^2$. Simplificando el primer término con la identidad suma por diferencia nos queda $g^{(\mathcal{D})}(x) = (x_2 + x_1)(x - x_1) + x_1^2$.

Sabemos por definición que $\bar{g}(x) = \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(x)]$. Como la distribución de los puntos es uniforme, su función de densidad es $f(x) = \frac{1}{\lambda(\mathcal{D})}$, con λ la medida del conjunto (que en este caso obviamente es 4). Además, x_1 y x_2 son independientes. Por tanto la esperanza quedaría:

$$\mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(x)] = \mathbb{E}_{\mathcal{D}} [(x_2 + x_1)(x - x_1) + x_1^2] = \frac{1}{\lambda(\mathcal{D})} \int_{\mathcal{D}} g^{(\mathcal{D})}(x) = \frac{1}{4} \int_{-1}^1 \int_{-1}^1$$

$$\begin{aligned} \bar{g}(x) = \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(x)] &= \mathbb{E}_{\mathcal{D}} [(x_2 + x_1)(x - x_1) + x_1^2] = \mathbb{E}_{\mathcal{D}} [x(x_1 + x_2) - x_1x_2] \\ &= x\mathbb{E}_{\mathcal{D}} [x_1 + x_2] - \mathbb{E}_{\mathcal{D}} [x_1x_2] \\ &= x(\mathbb{E}_{\mathcal{D}} [x_1] + \mathbb{E}_{\mathcal{D}} [x_2]) - \mathbb{E}_{\mathcal{D}} [x_1] \mathbb{E}_{\mathcal{D}} [x_2] \end{aligned}$$

Como el espacio de medida es simétrico (es un cuadrado), es claro que

$$\mathbb{E}_{\mathcal{D}} [x_1] = \mathbb{E}_{\mathcal{D}} [x_2]$$

Por tanto, veamos cuánto vale $\mathbb{E}_{\mathcal{D}} [x_1]$. Es claro que vale 0 ya que

$$\mathbb{E}_{\mathcal{D}} [x_1] = \frac{1}{4} \int_{-1}^1 \int_{-1}^1 x_1 dx_1 dx_2 = \frac{1}{4} \int_{-1}^1 \frac{x_1^2}{2} \Big|_{-1}^1 dx_2 = 0$$

Por tanto, sustituyendo todas las esperanzas en la fórmula de arriba, nos queda:

$$\bar{g}(x) = 0$$

Ahora calculemos E_{out} , *bias* y *var*. Aplicamos sus respectivas definiciones:

$$\begin{aligned} bias(x) &= (\bar{g}(x) - f(x))^2 = f(x)^2 = x^4 \\ var(x) &= \mathbb{E}_{\mathcal{D}} [(g^{(\mathcal{D})}(x) - \bar{g}(x))^2] = \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(x)^2] \end{aligned}$$

Esta varianza se calcula utilizando el módulo de Python *Sympy*, que es un completo paquete de cálculo simbólico.

```

1  import sympy
2
3  sympy.init_session()
4  x1, x2 = symbols('x1 x2')
5
6  print(integrate( ((x1+x2)*x - x1*x2) ** 2, (x1,-1,1),(x2,-1,1) ))
```

Nos devuelve

$$\frac{8x^2}{3} + \frac{4}{9}$$

Como queremos la esperanza en el conjunto y no solamente la integral, debemos multiplicarlo por la función densidad (la inversa de la medida del conjunto).

$$var(x) = \mathbb{E}_{\mathcal{D}} [(g^{(\mathcal{D})}(x) - \bar{g}(x))^2] = \mathbb{E}_{\mathcal{D}} [g^{(\mathcal{D})}(x)^2] = \frac{1}{4} \left(\frac{8x^2}{3} + \frac{4}{9} \right) = \frac{2x^2}{3} + \frac{1}{9}$$

Ahora sólo nos queda E_{out} . La esperanza en x se calcula igual que en \mathcal{D} , pero en este caso la medida del conjunto $[-1, 1]$ es 2.

$$E_{out} = \mathbb{E}_{\mathcal{D}} [E_{out}(g^{(\mathcal{D})})] = \mathbb{E}_x [bias(x) + var(x)] = \underbrace{\mathbb{E}_x [bias(x)]}_{=\frac{1}{5}} + \underbrace{\mathbb{E}_x [var(x)]}_{=\frac{1}{3}} = \frac{8}{15}$$

BONUS. Considere el enunciado del ejercicio 2 de la sección ERROR Y RUIDO de la relación de apoyo.

- a) Si su algoritmo busca la hipótesis h que minimiza la suma de los valores absolutos de los errores de la muestra,

$$E_{in}(h) = \sum_{n=1}^N |h - y_n|$$

entonces mostrar que la estimación será la mediana de la muestra, h_{med} (cualquier valor que deje la mitad de la muestra a su derecha y la mitad a su izquierda)

- b) Suponga que y_n es modificado como $y_n + \epsilon$, donde $\epsilon \rightarrow \infty$. Obviamente el valor de y_n se convierte en un punto muy alejado de su valor original. ¿Cómo afecta esto a los estimadores dados por h_{mean} y h_{med} ?

Solución. Colocamos h en $y_1 \leq \dots \leq y_k \leq h \leq y_{k+1} \leq \dots \leq y_N$. Descomponemos la sumatoria gracias a que todos los valores están ordenados:

$$\begin{aligned} E_{in}(h) &= \sum_{n=1}^N |h - y_n| = \sum_{n=1}^k (h - y_n) + \sum_{n=k+1}^N (y_n - h) = \sum_{n=1}^k h - \sum_{n=1}^k y_n + \sum_{n=k+1}^N y_n - \sum_{n=k+1}^N h \\ E_{in}(h) &= kh - \sum_{n=1}^k y_n + \sum_{n=k+1}^N y_n - (N - k)h \end{aligned}$$

Ahora que tenemos los signos arreglados, procedemos a buscar un mínimo, derivando e igualando a cero.

$$E'_{in}(h) = k - (N - k) = 2k - N = 0 \Leftrightarrow k = \frac{N}{2}$$

Hemos determinado que para que h sea un mínimo, debe estar entre y_k y y_{k+1} con $k = \frac{N}{2}$, que hace que esté en medio de todos los datos. Por tanto, este valor es la mediana. /*Añadir algo para argumentar el caso N impar*/

BONUS. Considere el ejercicio 12.

- a) Describir un experimento que podamos ejecutar para determinar (numéricamente) $g(x)$, E_{out} , $bias$ y var .
- b) Ejecutar el experimento y dar los resultados. Comparar E_{out} con $bias+var$. Dibujar en unos mismos ejes $g(x)$, E_{out} y $f(x)$.

Solución.